

## PATENT ABSTRACTS OF JAPAN

(11)Publication number : 2001-222517

(43)Date of publication of application : 17.08.2001

(51)Int.Cl.

G06F 15/177

G06F 12/08

G06F 13/16

G06F 15/16

(21)Application number : 2000-111165

(71)Applicant : CHO SEITAI  
ZEN SHUSHOKU  
KIN MEICHU

(22)Date of filing : 12.04.2000

(72)Inventor : CHO SEITAI  
ZEN SHUSHOKU  
KIN MEICHU

(30)Priority

Priority number : 2000 200006401

Priority date : 11.02.2000

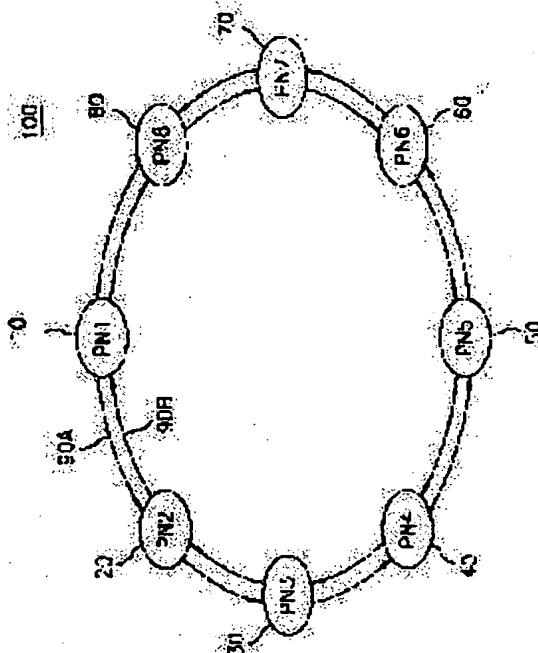
Priority country : KR

## (54) DISTRIBUTED SHARED MEMORY MULTIPLEX PROCESSOR SYSTEM HAVING DIRECTION-SEPARATED DUPLEX RING STRUCTURE

## (57)Abstract:

**PROBLEM TO BE SOLVED:** To provide a distributed shared memory multiplex processor system having direction-separated duplex ring structure while using a snooping system.

**SOLUTION:** There are plural processor nodes arrayed in the shape of ring, any one of plural processor nodes generates a request signal to one data block and the remaining processor nodes snoop their own internal elements. Thus, there are plural processor nodes for supplying a data block from any one of the remaining processor nodes and there is a direction-separated duplex ring structure for supplying two opposite routes along with first and second ring paths while including the first and second ring paths by snooping the internal element of the remaining processor nodes. Then, this system is provided with the direction-separated duplex ring structure, to which the plural processor nodes are connected through the first and second routes, for performing the multi-address communication of the request signal through any one of routes to each of the remaining processor nodes and performing the single communication of the data block through any one of routes to the processor node, which generates the request signal.



## LEGAL STATUS

[Date of request for examination]

04.04.2002

[Date of sending the examiner's decision of rejection]

06.06.2006

[Kind of final disposal of application other than the examiner's decision of rejection or application converted registration]

[Date of final disposal for application]

[Patent number]

[Date of registration]

[Number of appeal against examiner's decision of rejection]

[Date of requesting appeal against examiner's decision of rejection]

[Date of extinction of right]

(19)日本国特許庁 (J P)

(12) 公 開 特 許 公 報 (A)

(11)特許出願公開番号

特開2001-222517

(P2001-222517A)

(43)公開日 平成13年8月17日(2001.8.17)

(51)Int.Cl.	識別記号	F I	テーマコード(参考)
G 0 6 F 15/177 12/08	6 8 2	G 0 6 F 15/177 12/08	6 8 2 A 5 B 0 0 5 E 5 B 0 4 5
	3 1 0		3 1 0 B 5 B 0 6 0
13/16	5 1 0	13/16	5 1 0 B
15/16	6 4 0	15/16	6 4 0 M
審査請求 未請求 請求項の数13 O L (全 18 頁)			

(21)出願番号 特願2000-111165(P2000-111165)

(22)出願日 平成12年4月12日(2000.4.12)

(31)優先権主張番号 2 0 0 0 - 6 4 0 1

(32)優先日 平成12年2月11日(2000.2.11)

(33)優先権主張国 韓国 (K R)

(71)出願人 598111009

張 星泰

大韓民国、ソウル特別市城北区貞陵1洞  
1015番地 京南アパートメント103-1701

(71)出願人 598111010

全 洲植

大韓民国、ソウル特別市江南区道谷洞星宇  
アパートメント1-103

(74)代理人 100058479

弁理士 鈴江 武彦 (外5名)

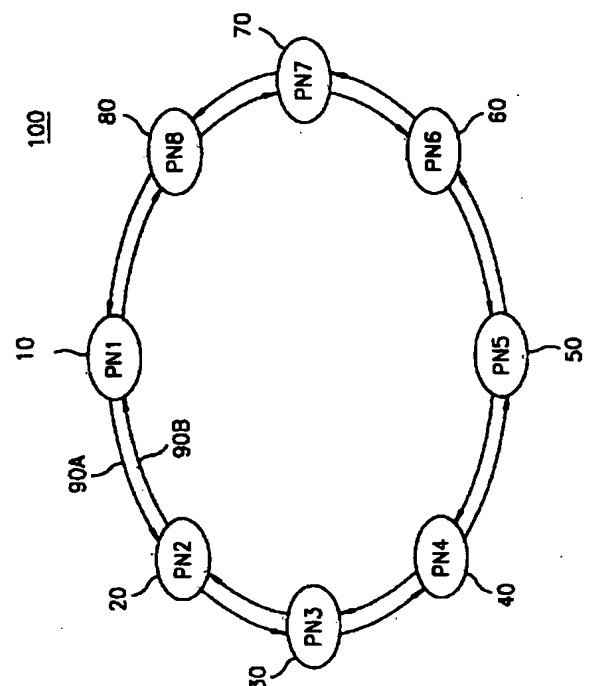
最終頁に続く

(54)【発明の名称】 方向分離二重リング構造を有する分散共有メモリ多重プロセッサシステム

(57)【要約】

【課題】 スヌーピング方式を用い方向分離二重リング構造を有する分散共有メモリ多重プロセッサシステムを提供する。

【解決手段】 リング状で配列されている複数のプロセッサノードがあり、複数のプロセッサノードのいずれかが1つのデータブロックに対する要求信号を発生し、残余プロセッサノードが自分の内部要素をスヌーピングすることにより、残余プロセッサノードのいずれかがデータブロックを供給する複数のプロセッサノードと、第1、及び第2リングパスを含めて第1、及び第2リングパスに沿って2つの反対経路を供給する方向分離二重リング構造があって、第1及び第2経路を介して、複数のプロセッサノードが接続され、要求信号が経路のいずれかを介して残余プロセッサノードの各々に同報通信され、データブロックは経路のいずれかを介して要求信号を生成したプロセッサノードに単一通信される方向分離二重リング構造を含む。



## 【特許請求の範囲】

【請求項 1】 分散共有メモリ多重プロセッサシステムであって、

リング状で配列されている複数のプロセッサノードがあり、前記複数のプロセッサノード中のいずれかが 1 つのデータブロックに対する要求信号を発生し、残余プロセッサノードが自分の内部要素をスヌーピングすることにより、前記残余プロセッサノード中のいずれが前記データブロックを供給する前記複数のプロセッサノードと、第 1 及び第 2 リングバスを含み、前記第 1 及び第 2 リングバスに沿って 2 つの反対経路を供給する方向分離二重リング構造があり、前記第 1 及び第 2 経路を介して前記複数のプロセッサノードが接続され、前記要求信号が前記経路中のいずれかを介して前記残余プロセッサノードの各々に同報通信され、前記データブロックは前記経路中のいずれかを介して前記要求信号を生成したプロセッサノードに単一通信される前記方向分離二重リング構造とを備えることを特徴とする分散共有メモリ多重プロセッサシステム。

【請求項 2】 前記要求信号は、前記データブロックが前記複数のプロセッサノードの内部要素の偶数または奇数メモリブロックアドレス中のどこに格納されているかによって、前記第 1 リングバスまたは第 2 リングバスを介して伝送されることを特徴とする請求項 1 に記載の分散共有メモリ多重プロセッサシステム。

【請求項 3】 前記要求信号を生成したプロセッサノードが前記要求信号に対応する対応信号を受信する前に前記残余プロセッサノード中の少なくとも一つから前記データブロックとは異なるデータブロックに対する要求信号を受信すると、前記プロセッサノードは、まず前記残余プロセッサノード中の少なくとも一つから伝達された要求信号に対応する動作を行った後、前記応答信号に対応する動作を処理することを特徴とする請求項 1 に記載の分散共有メモリ多重プロセッサシステム。

【請求項 4】 前記複数のプロセッサノードの各々は前記データブロックに対する前記要求信号を発生する複数のプロセッサモジュールと、前記複数のプロセッサモジュールの各々により共有されるデータブロックを格納するローカル共有メモリと、前記要求信号に反応して、前記要求信号に応じる前記データブロックが前記ローカル共有メモリに有効な状態で格納されているかどうかを調査し、前記データブロックが前記ローカル共有メモリ内に有効な状態で格納されている場合は、前記データブロックを前記複数のプロセッサモジュールに伝達し、前記データブロックが前記ローカル共有メモリ内に有効な状態で格納されていない場合は、前記要求信号を前記複数のプロセッサノード中の隣合う次のプロセッサノードに供給するノード制御器と、前記ノード制御器と前記方向分離二重リング構造をインタフェースするリンク制御器と、

前記複数のプロセッサモジュールとローカル共有メモリとノード制御器とを相互接続する相互接続網とを含むことを特徴とする請求項 1 に記載の分散共有メモリ多重プロセッサシステム。

【請求項 5】 前記複数のプロセッサノードの各々は、前記データブロックに対する前記要求信号を発生する複数のプロセッサモジュールと、遠隔キャッシュと、前記要求信号に反応して、前記要求信号に対応する前記データブロックが前記遠隔キャッシュに有効な状態で格納されているかどうかを調査し、前記データブロックが前記遠隔キャッシュ内に有効な状態で格納されている場合は、前記データブロックを前記複数のプロセッサモジュールに伝達し、前記データブロックが前記遠隔キャッシュ内に有効な状態で格納されていない場合は、前記要求信号を前記複数のプロセッサノードの中の隣合う次のプロセッサノードに供給するノード制御器と、前記ノード制御器と前記方向分離二重リング構造をインタフェースするリンク制御器と、

前記複数のプロセッサモジュールとローカル共有メモリとノード制御器とを相互接続する相互接続網とを含むことを特徴とする請求項 1 に記載の分散共有メモリ多重プロセッサシステム。

【請求項 6】 前記複数のプロセッサノードの各々は、前記データブロックに対する前記要求信号を発生する複数のプロセッサモジュールと、前記複数のプロセッサモジュールの各々により共有されるデータブロックを格納するローカル共有メモリと、遠隔キャッシュと、

前記要求信号に反応して、前記要求信号に対応する前記データブロックが前記遠隔キャッシュまたは前記ローカル共有メモリに有効な状態で格納されているかどうかを調査し、前記データブロックが前記遠隔キャッシュまたは前記ローカル共有メモリ内に有効な状態で格納されている場合は、前記データブロックを前記複数のプロセッサモジュールに伝達し、前記データブロックが前記遠隔キャッシュまたは前記ローカル共有メモリ内に有効な状態で格納されていない場合は、前記要求信号を前記複数のプロセッサノード中の隣合う次のプロセッサノードに供給するノード制御器と、

前記ノード制御器と前記方向分離二重リング構造をインタフェースするリンク制御器と、前記複数のプロセッサモジュールとローカル共有メモリとノード制御器とを相互接続する相互接続網とを含むことを特徴とする請求項 1 に記載の分散共有メモリ多重プロセッサシステム。

【請求項 7】 前記複数のプロセッサノードの各々は、前記ローカル共有メモリに格納されたデータブロックに対する状態情報を含むメモリディレクトリをさらに含む請求項 4 または請求項 6 に記載の分散共有メモリ多重プロセッサシステム。

ロセッサシステム。

【請求項8】 前記メモリディレクトリは2つあるいはそれ以上の独立的なメモリディレクトリを有することにより、前記プロセッサモジュール及び前記残余プロセッサノードからの前記ローカル共有メモリに対するアクセス要求を並列に処理することを特徴とする請求項7に記載の分散共有メモリ多重プロセッサシステム。

【請求項9】 前記遠隔キャッシュは、前記データブロックの内容を含む遠隔データキャッシュと、前記遠隔データキャッシュに格納された前記データブロックのタグアドレス及び状態を格納する遠隔タグキャッシュとを含むことを特徴とする請求項5または請求項6に記載の分散共有メモリ多重プロセッサシステム。

【請求項10】 前記遠隔タグキャッシュが、2つあるいはそれ以上の独立的な遠隔タグキャッシュを有することにより、前記プロセッサモジュール及び前記残余プロセッサノードからの前記遠隔キャッシュに対するアクセス要求を並列に処理することを特徴とする請求項9に記載の分散共有メモリ多重プロセッサシステム。

【請求項11】 前記リンク制御器は、前記ノード制御器に前記第1リンクバス及び第2リンクバスをそれぞれ接続する第1リンク制御器と第2リンク制御器とを含むことを特徴とする請求項4ないし請求項6のいずれかに記載の分散共有メモリ多重プロセッサシステム。

【請求項12】 前記第1リンク制御器と前期第2リンク制御器の各々は、前記要求信号と前記データブロックとを含むバケットを生成することにより、各々の制御器に接続された対応リングバスを介して前記残余プロセッサノードに前記バケットを伝送し、前記残余プロセッサノードから前記対応リングバスを介して供給された要求またはデータブロックを前記ノード制御器に選択的に供給することを特徴とする請求項11に記載の分散共有メモリ多重プロセッサシステム。

【請求項13】 前記プロセッサノードの各々は、前記ノード制御器と前記第1リンク制御器及び前記第2リンク制御器を接続するリンクバスをさらに含むことを特徴とする請求項12に記載の分散共有メモリ多重プロセッサシステム。

【発明の詳細な説明】

【0001】

【発明の属する技術分野】 本発明は、分散共有メモリ多重プロセッサシステムに関し、より詳しくは方向分離二重リング構造を備えた分散共有メモリ多重プロセッサシステムに関する。

【0002】

【従来の技術】 一般的に、多重プロセッサシステムは明示的なメッセージ伝送(message-passing)を用いてプロセッサ間の通信を実現する分散メモリ構造(distributed-memory architecture)と、メモリを共有して単一システ

ムイメージを供給する共有メモリ構造(shared-memory architecture)に分けることができるが、現在共有メモリ多重プロセッサシステムの相互連結網の中で、もっとも大衆的な技術は共有バスである。共有バスは実現上の複雑度が低く、且つ低費用であるということで広く使用されているが、性能が速やかに向上しているプロセッサの速度に追いつくことができないという短所を有する。また、バスの物理的な特性による拡張性や、バス使用量の増加によるバス帯域幅(bandwidth)に不具合がある。

10 【0003】このようなバス構造の限界を克服するためにいろいろ試してきた結果、単方向地点間リンクを用いたIEEE SCIが標準として確定された。最大四つのプロセッサがスヌーピング方式のバスによりUMA(Uniform Memory Access)の形で設けられたプロセッサノードをSCIリンクを用いて単方向リング構造で接続し、ディレクトリ方式のキャッシュプロトコルを利用して実現された常用化システムが開示されている。(Tom LovettとRussell Clappの“STING: A CC-NUMA Computer System for the Commercial Marketplace” (In Proceedings of the 23th International Symposium on Computer Architecture, pp.308-317, May 1996を参照) 20 上述したシステムをさらに改善したもので、図1及び図2に示すように、最大四つのプロセッサがスヌーピング方式のバスによりUMAの形で設けられたプロセッサノードをSCIリンクを用いて単方向リング構造で接続したシステムをスヌーピング方式のキャッシュプロトコルを用いて具現した多重プロセッサシステムが張星泰と全洲植と鄭盛宇の“PANDA: Ring-Based Multiprocessor System using New Snooping Protocol” (In The Proceeding of ICPADS 1998, pp.10-17, Dec.1998) (1998年8月7日、日本国出願された特許出願第224423/98号を参照)に開示されている。

【0004】しかし、プロセッサとローカルバスのクラック速度が向上しつつあることに従って、このような高性能のプロセッサとローカルバスを採択した前記分散共有メモリ多重プロセッサシステムは単一地点間リンク帯域幅及びシステム拡張が必要となった。それにより、既存システムで単一地点間リンク帯域幅を拡張するために、単に2倍の帯域幅を有するリンクを使用する方法もあるが、現在2倍増加された帯域幅を有する新しいリンクの開発、及びこれを短時間内に常用化された製品に適用するのは実際には難しい。

40 【0005】

【発明が解決しようとする課題】 従って、本発明の主な目的は、スヌーピング方式を用いてシステム内のプロセッサノード間のキャッシュ一貫性を維持しながら、且つ性能を向上するために方向分離二重リング構造を有する分散共有メモリ多重プロセッサシステムを提供することにある。

50 【0006】

【課題を解決するための手段】上記の目的を達成するために、本発明によれば、リング状で配列されている複数のプロセッサノードがあり、前記複数のプロセッサノード中のいずれかが1つのデータブロックに対する要求信号を発生し、残余プロセッサノードが自分の内部要素をスヌーピングすることにより、前記残余プロセッサノード中のいずれかが前記データブロックを供給する前記複数のプロセッサノードと、第1及び第2リングバスを含み、前記第1及び第2リングバスに沿って2つの反対経路を供給する方向分離二重リング構造があり、前記第1及び第2経路を介して前記複数のプロセッサノードが接続され、前記要求信号が前記経路中のいずれかを介して前記残余プロセッサノードの各々に同報通信され、前記データブロックは前記経路中のいずれかを介して前記要求信号を生成したプロセッサノードに単一通信される前記方向分離二重リング構造とを備える分散共有メモリ多重プロセッサシステムが提供される。

#### 【0007】

【発明の実施の形態】以下、本発明は添付の図面を参照して、次のように詳細に説明する。

【0008】図3を参照すると、スヌーピング方式を支援する地点間(point-to-point)方向分離二重リング構造、例えば、90A及び90Bに基づく分散共有メモリ多重プロセッサシステム100を示し、方向分離二重リングバス90A及び90Bは地点間リンクを用いて実現し、地点間リンクの各々は複数の信号を伝送することができる光ファイバ、同軸ケーブルまたは光連結部を用いて実現することができる。本発明の好適な実施例によると、分散共有メモリ多重プロセッサシステム100は8個のプロセッサノード、即ちPN1~PN8(10~80)を含む。PN1(10)からPN8(80)に達する各々のプロセッサノードは、スヌーピングを支援する地点間方向分離二重リング90A及び90Bを介して接続されている。

【0009】ここで、「方向分離二重リング」とは、進行方向が互いに反対である2つの分離されたリングバスを有し、プロセッサノードからのデータ要求信号またはその要求信号に対応するデータブロックを一番近い経路を有する第1方向、または前記第1方向と反対である第2方向に伝達するように構成されたリングバス構造を意味する。図3に示すように、方向分離二重リング構造は第1リングバス90Aと第2リングバス90Bを含み、要求信号及び検索されたデータブロックはデータブロックの属性、例えば、あるプロセッサノードにより要求されたデータを含むデータブロックが偶数メモリブロックアドレスと奇数メモリブロックアドレスの中のどこに格納されているかにより、第1リングバス90Aまたは第2リングバス90Bを介して伝送される。

【0010】図4には、互いに同一の構成を有する多数個のプロセッサノードPN1~PN8(10~80)中のいずれかのプロセッサノード、例えば、PN1(10)の構成がより詳細

に示されている。図4に示すように、PN1(10)はキャッシュを組み込んでいる多数個のプロセッサモジュールと、入出力(I/O)ブリッジ216と、ローカルシステムバス218と、ローカル共有メモリ220と、ノード制御器222と、遠隔キャッシュ224と、メモリディレクトリ226と、リンクバス228と、リングインタフェース230とを含む。

【0011】説明の便宜上、単なる2つのプロセッサモジュール、即ち第1プロセッサモジュール212と、第2プロセッサモジュール214が図4に示されており、プロセッサモジュール212及び214と、ローカル共有メモリ220と、I/Oブリッジ216と、ノード制御器222は、ローカルシステムバス218を通じて互いに接続される。

【0012】リングインタフェース230は、2つのリンク制御器230A及び230Bを含み、このリンク制御器を介して各々のプロセッサノードは方向分離二重リング構造、即ちリングバス90A及び90Bにそれぞれ接続される。この実施例で、リングインタフェース230は、リンクバス228を介してノード制御器222に接続される。

【0013】ノード制御器222は、プロセッサモジュール212及び214のいずれかからの要求信号に対応するデータブロックが遠隔キャッシュ224またはローカル共有メモリ220に有効な状態で格納されているかどうかを検索する。検索の結果、データブロックが遠隔キャッシュ224に有効な状態で格納されている場合は、ノード制御器222は遠隔キャッシュ224に格納された該当データブロックを要求信号を発生させるプロセッサモジュールに供給するが、データブロックがローカル共有メモリ220に有効な状態で格納されている場合は、ローカル共有メモリ220が要求信号を発生させたプロセッサモジュールにデータブロックを供給する。

【0014】遠隔キャッシュ224やローカル共有メモリ220ともにそのデータブロックが有効な状態で格納されておらず、データブロックが奇数ブロックアドレス、即ち奇数ブロックである場合は、そのデータブロックに対する要求信号を第1リンク制御器230Aを介して第1リングバス90Aに伝送する。一方、データブロックが偶数ブロックアドレス、即ち偶数ブロックである場合は、そのデータブロックに対する要求信号を第2リンク制御器230Bを介して第2リングバス90Bに伝送する。このようにノード制御器222はリングインタフェース230、即ち第1リンク制御器230A及び第2リンク制御器230Bを介してデータブロックに対する要求信号を他のプロセッサノードPN2~PN8(20~80)に供給する。

【0015】続いて、他のプロセッサノードPN2~PN8(20~80)のいずれかからデータブロックに対する要求信号が第1リンク制御器230Aまたは第2リンク制御器230Bを介して受信されると、ノード制御器222はその要求信号に対応するデータブロックが自分の遠隔キャッシュ224やローカル共有メモリ220に有効な状態で格納されているかどうかを検索する。検索の結果、データブロックが遠隔キャ

ッシュ224またはローカル共有メモリ220に有効な状態で格納されている場合は、ノード制御器222はローカルシステムバス218を介して遠隔キャッシュ224またはローカル共有メモリ220からデータブロックを受信して、そのデータブロックを要求したプロセッサノードに一番近い経路を有する第1リングバス90Aまたは第2リングバス90Bを介してそのデータブロックを伝送する。

【0016】図4に示すように、各々のプロセッサノードPN1~PN8(10~80)は、ローカル共有メモリ220に格納されたデータブロックに対する状態情報を格納するメモリディレクトリ226をさらに含み、ノード制御器222が直接メモリディレクトリ226をアクセスする。従って、メモリディレクトリ226によりノード制御器222はプロセッサモジュール212及び214のいずれかから要求されたデータブロックがローカル共有メモリ220にどの状態で格納されているかを効果的に検索し、他のプロセッサノードPN2~PN8(20~80)のいずれかから要求されたデータブロックが自分のローカル共有メモリ220にどの状態で格納されているかを効果的に検索することができるようになる。さらに好ましくは、メモリディレクトリ226は、図4に示すように、独立的な2つのメモリディレクトリ226A及び226Bより構成される。

【0017】第1メモリディレクトリ226Aは、ローカルシステムバス218を介して伝達されるプロセッサモジュール212及び214からのローカル共有メモリアクセス要求に反応し、第2メモリディレクトリ226Bは、リンクバス228を介してノード制御器222に接続されたリングインタフェース230を介して他のプロセッサノードPN2~PN8(20~80)から伝達された遠隔共有メモリアクセス要求に反応する。このような方式を通じて、ローカルメモリ220に対するアクセス要求が並列に行われるようにすることができる。

【0018】第1リンク制御器230Aと第2リンク制御器230Bは、PN1(10)を方向分離二重リングバス90A及び90Bに接続するデータ経路を供給し、パケット伝送に必要な全体的なデータの流れを制御する。第1リンク制御器230Aと第2リンク制御器230Bは、ノード制御器222からの要求信号やデータブロックを有するパケットを構成して第1リングバス90Aまたは第2リングバス90Bを介して他のプロセッサノードPN2~PN8(20~80)に伝送し、第1リングバス90Aまたは第2リングバス90Bを介して他のプロセッサノードPN2~PN8(20~80)から伝送されてくる要求信号やデータブロックを選別してノード制御器222に伝送する。また、リンクインタフェース230は伝送される要求信号が放送パケットである場合は、その放送パケットをスヌーピングのためにノード制御器222に伝送するだけでなく、伝送された放送パケットを次のプロセッサノードPN2(20)またはPN8(80)にバイパス(bypass)する。より具体的に、第1リンク制御器230Aは第1リングバス90Aを介して隣り合うプロセッサノードPN8(80)から伝達され

る放送パケットを第1リングバス90Aを介して隣り合うさらに他のプロセッサノードPN2(20)に伝達し、第2リンク制御器230Bは第2リングバス90Bを介して隣り合うプロセッサノードPN2(20)から伝送される放送パケットを第2リングバス90Bを介して隣り合うさらに他のプロセッサノードPN8(80)に伝達する。

【0019】一方、遠隔キャッシュ224は自分を除く他のプロセッサノードPN2~PN8(20~80)のローカル共有メモリ(以下、遠隔共有メモリと称する)に格納されたデータブロックのみをキャッシュする。ローカルシステムバス218に接続されたプロセッサモジュール212及び214のいずれかが他のプロセッサノードPN2~PN8(20~80)のいずれかの遠隔共有メモリに格納されたデータブロックを要求する場合、そのデータブロックは遠隔キャッシュ224に割当てられ、ローカル共有メモリ220に格納されたデータブロックはキャッシングされない。即ち、上述のように、プロセッサノードPN1(10)の遠隔キャッシュ224は他のプロセッサノードPN2~PN8(20~80)の遠隔共有メモリに格納されたデータブロックのみをキャッシングすることにより、遠隔メモリアクセス時間を減らすことができる。

【0020】遠隔キャッシュ224はプロセッサノードPN1(10)内のプロセッサモジュール212及び214におけるキャッシュや、他のプロセッサノードPN2~PN8(20~80)内の遠隔共有メモリに対してMLI性質(Multi-Level Inclusion Property)を満たすため、他のプロセッサノードPN2~PN8(20~80)からの遠隔共有メモリ参照要求に対するスヌーピング・フィルタリング(Snoop filtering)機能を行うことができる。ここで、MLI性質は、下位階層、即ちローカルキャッシュに有効な状態で格納されたデータブロックは上位階層、即ち遠隔キャッシュにも常に有効な状態で格納されていなければならない性質を意味する。このようなMLI特性を保障するため、上位階層のキャッシュに格納されたデータブロックが置換(replacement)えられる場合、該当データブロックは下位階層のどのキャッシュにも有効な状態で存在してはならない。

【0021】従って、遠隔キャッシュ224はプロセッサノードPN1(10)のプロセッサモジュール212及び214内のキャッシュに有効な状態で格納された遠隔データブロックを格納ようになる。他のプロセッサノードPN2~PN8(20~80)からの遠隔共有メモリ参照要求信号に反応するデータブロックが遠隔キャッシュ224に有効な状態で格納されていない場合は、ローカルシステムバス218で該当データブロックに対する要求を伝送する必要がないスヌーピング・フィルタリング機能を行う。

【0022】このとき、好ましくは遠隔キャッシュ224はデータブロックの内容を格納する遠隔データキャッシュ224-1と、データブロックの状態及びアドレスの一部を格納する遠隔タグキャッシュ224-2を含むことにより、遠隔キャッシュ224に格納されたデータブロックの

状態を更新するか、または必要な場合、該当データブロックを供給し易くする。さらに好ましくは、遠隔タグキャッシュ224-2は遠隔データブロックのアドレスや状態を格納する2つの遠隔タグキャッシュ、即ち第1遠隔タグキャッシュ224-2Aと第2遠隔タグキャッシュ224-2Bを含む。第1遠隔タグキャッシュ224-2Aは、プロセッサモジュール212及び214のいずれかによる遠隔キャッシュアクセス要求に反応し、第2遠隔タグキャッシュ224-2Bは他のプロセッサノードPN2~PN8(20~80)のいずれかによる遠隔キャッシュアクセス要求に反応する。このような方式で、遠隔キャッシュ224に対するアクセス要求を並列に処理することができる。

【0023】上述したように構成された本発明の方向分離二重リング構造を有する分散共有メモリ多重プロセッサシステムの動作は図5及び図6を参照して次のように詳しく説明する。

【0024】遠隔キャッシュ224に格納されたデータブロックは次の四つの状態、即ち「更新(Modified)」、「更新-共有(Modified-Shared)」、「共有(Shared)」、「無効(Invalid)」状態のいずれかで表れる。四つの状態はそれぞれ次のように述べられる。

\*更新: 該当データブロックが有効で更新されており、唯一に有効なコピーである状態

\*更新-共有: 該当データブロックが有効で更新されており、他の遠隔キャッシュが該当データブロックを共有することができる。

\*共有: 該当データブロックが有効であり、他の遠隔キャッシュが該当データブロックを共有することができる。

\*無効: 該当データブロックが有効ではない。

【0025】また、本発明による多重プロセッサシステムで、第1及び第2メモリディレクトリ226A及び226Bは三つの状態、即ちC(clean)、S(share)、G(gone)のいずれかを維持することにより、ローカルシステムバス218を介してローカル共有メモリアクセス要求に反応するキャッシュコピーレントのトラフィック量を最小化し、方向分離二重リングバス90A及び90Bへの不必要なトランザクションを減らし、ローカルシステムバス219からの要求を効果的に処理し、方向分離二重リングバス90A及び90Bによるスヌーピング要求に反応するスヌーピング結果を発生する。この三つの状態に対する詳細は次に述べる。

\*C(Clean): 該当データブロックが他のプロセッサノードのどの遠隔キャッシュにも有効な状態で格納されていない。

\*S(Shared): 該当データブロックが有効であり、他のプロセッサノードのいずれかの遠隔キャッシュに更新されない有効な状態で格納され得る。

\*G(Gone): 該当データブロックが有効ではなく、他のプロセッサノードのいずれかの遠隔キャッシュに更新された有効な状態で格納されている。

【0026】一方、各プロセッサノードPN1~PN8(10~80)を順次接続する方向分離二重リングバス90A及び90B上の全ての通信はパケットを介して構成され、パケットは要求パケット、応答パケット、認識パケットに分類することができる。要求パケットは、方向分離二重リングバス90A及び90Bへのトランゼクションを必要とするプロセッサノードPN1~PN8(10~80)のいずれかにより発生されるパケットで同報通信パケット(broadcast packet)と単一通信パケット(unicast packet)に区分できる。この中で同報通信パケットのみが他のプロセッサノードによりスヌーピングされる。

【0027】応答パケットは要求パケットを受信したプロセッサノードにより常に単一通信される。認識パケットは、応答パケットを受信したプロセッサノードにより発生すると、応答パケットを伝送したプロセッサノードに単一通信される。応答パケットを単一通信したプロセッサノードは応答パケットに対応する認識パケットが受信されるまで応答パケットを維持する。本発明の他の実施例によると、応答パケットを単一通信したプロセッサノードが応答パケットに対応する認識パケットを受信する前に同一のデータブロックに対する他の要求パケットを他のプロセッサノードから受信した場合は、必要に応じて、そのプロセッサノードで要求パケットを再伝送することを要求することができる。

【0028】図5及び図6で、分散共有メモリ多重プロセッサシステムの例示的な動作が示されている。

【0029】第1番目の場合は、プロセッサノードPN1(10)内の第1プロセッサモジュール212がデータブロックに対する読込み要求を発生させる場合で、要求パケットはRQ12である。該当データブロックが遠隔共有メモリに該当し、PN1(10)の遠隔キャッシュ224に有効な状態で格納されていない場合は、PN1(10)はリングバス90Aまたは90Bを介して他のプロセッサノードPN2~PN8(20~80)にRQ12を同報通信する。また、該当データブロックがローカル共有メモリ220には該当するが、ローカル共有メモリ220に有効な状態で格納されていない場合は、PN1(10)はリングバス90Aまたは90Bを介して他のプロセッサノードPN2~PN8(20~80)にRQ12を通報通信する。この場合、該当データブロックが奇数ブロックであれば、RQ12は第1リンク制御器230A及び第1リンクバス90Aを介して伝送され、該当データブロックが偶数ブロックであれば、RQ12は第2リンク制御器230Bと第2リンクバス90Bを介して伝達される。

【0030】これにより、奇数ブロックに対するRQ12は、図5に示すように、第1リングバス90Aに沿って反時計方向にPN2(20)からPN8(80)に巡回する。一方、偶数ブロックに対するRQ12は、図6に示すように、第2リングバス90Bに沿って時計方向にPN8(80)からPN2(20)に巡回する。RQ12がリングバス90Aまたは90Bに沿って巡回する間、各々のプロセッサノードはRQ12に反応して内部の遠

隔キャッシュまたはメモリディレクトリを調査して該当データブロックがどの状態で格納されているかなどに対するスヌーピングを行うと同時に、そのRQ12を隣合う次のプロセッサノードにバイパスする。

【0031】例えば、RQ12がPN4(40)に供給されると、PN4(40)のノード制御器はPN4(40)内の遠隔キャッシュとメモリディレクトリをスヌーピングする。その結果、該当データブロックがPN4(40)の遠隔キャッシュに「更新」または「更新-共有」状態で格納されている場合は、PN4(40)のノード制御器はRQ12に応答する責任を有すると判断する。この場合、該当データブロックはローカル共有メモリに有効な状態で格納しているプロセッサノードは存在しない。

【0032】これにより、PN4(40)のノード制御器は該当データブロックを含む応答パケットRSP42をPN1に一番近い経路を有する第1リングバス90Aまたは第2リングバス90Bに沿ってPN1(10)に単一通信する。また、PN4(40)のノード制御器はPN4(40)の遠隔キャッシュの状態を「更新-共有」または「共有」状態のように更新されていない有効な状態に変更させる。

【0033】該当データブロックがPN4(40)のローカル共有メモリに有効な状態で格納されている場合は、PN4(40)のノード制御器はRQ12に対する応答の責任を取り、該当データブロックを含む応答パケットをPN1に一番近い経路を有する第1リングバス90Aまたは第2リングバス90Bを介してPN1(10)に伝送する。それで、図5及び図6で、PN4(40)はPN1(10)まで一番近い経路を有する第2リングバス90Bを介して応答パケットRSP42を伝送する。

【0034】このとき、前記読込み要求パケットRQ12は巡回を終了した後、PN1(10)により除去される。一方、RSP42を受信した後、プロセッサノードPN1(10)はプロセッサノードPN4(40)まで一番近い経路を有する第1リングバス90Aまたは第2リングバス90B、即ち図5及び図6では第1リングバス90Aを介して認識パケットACK14をPN4(40)に単一通信すると同時に、プロセッサノードPN1(10)のローカルシステムバス218を介して該当データブロックに対する読込み要求を発生させる第1プロセッサモジュール212にRSP42内の該当データブロックを伝送する。該当データブロックが遠隔共有メモリに該当すると、PN1(10)は遠隔データキャッシュ224-1に該当データブロックを有効な状態で格納すると同時に、該当データブロックに対応する遠隔タグキャッシュ224-2の状態を有効な状態に更新し、該当データブロックがローカル共有メモリ220に該当すると、PN1(10)はローカル共有メモリ220に該当データブロックを格納すると同時に、メモリディレクトリ226の状態を他のプロセッサノードが該当データブロックを共有していることを意味する状態、例えば「S」状態に更新する。

【0035】一方、第2番目の場合は、PN1(10)内の第1プロセッサモジュール212が書込み要求を発生する場合

であり、この場合にも図5及び図6はまだ有効であり、単純化のためにこの時の要求パケットもRQ12とする。PN1(10)が書込み要求に対する該当データブロックを遠隔キャッシュ224やローカル共有メモリ220のどこにも有効な状態で格納していないと、PN1(10)はリングバス90Aまたは90Bを介して他のプロセッサノードPN2~PN8(20~80)にRQ12を同報通信する。この場合、該当データブロックが奇数ブロックである場合は、RQ12は、図5に示すように、第1リングバス90Aを介して伝送される。一方、該当データブロックが偶数ブロックである場合は、RQ12は、図6に示すように、第2リングバス90Bを介して伝送される、書込み要求パケットRQ12がリングバス90Aまたは90Bを介して巡回する間、各々のプロセッサノードはRQ12に反応して内部の遠隔キャッシュまたはメモリディレクトリを調査して該当データブロックがどのような状態で格納されているかなどに対するスヌーピングを行うと同時に、RQ12を隣合う次のプロセッサノードにバイパスする。

【0036】例えば、RQ12がPN4(40)に供給されると、PN4(40)のノード制御器は内部の遠隔キャッシュとメモリディレクトリをスヌーピングする。その結果、該当データブロックがPN4(40)の遠隔キャッシュに更新された状態、例えば、「更新」または「更新-共有」状態で格納されているか(この場合に、該当ブロックをローカル共有メモリに有効な状態で格納しているプロセッサノードは存在しない)、または該当データブロックがローカル共有メモリに有効な状態で格納されていると、PN4(40)のノード制御器は自分がRQ12に対する応答の責任を有すると判断して要求したデータブロックを含む応答パケットRSP42をPN1(10)に一番近い経路を有する第1リングバス90Aまたは第2リングバス90Bを介して伝送する。図5及び図6を参照すると、PN4(40)は一番近い経路を有する第2リングバス90Bを介して応答パケットRSP42を伝送する。また、PN4(40)のノード制御器は該当データブロックを格納している遠隔キャッシュの状態を無効化された状態、例えば、「無効」状態にするか、またはメモリディレクトリの状態で無効化された状態、例えば、「G」状態に更新する。前記書込み要求パケットRQ12は、リングバス90Aまたは90Bを巡回した後、プロセッサノードPN1(10)により除去される。

【0037】一方、スヌーピング結果、該当データブロックが他のプロセッサノード、即ちPN2~PN3(20~30)とPN5~PN8(50~80)の遠隔キャッシュに更新されていない有効な状態、例えば、「共有」状態で格納されていることと判明されると、プロセッサノード、即ちPN2~PN3(20~30)とPN5~PN8(50~80)の遠隔キャッシュ状態は無効化された状態、例えば、「無効」状態に変更する。

【0038】PN1(10)はPN4(40)からRSP42を受信すると、PN4(40)に一番近い経路を有する第1リングバス90Aまたは第2リングバス90B、例えば、図5及び図6の場合は



第1リングバス90Aを介してPN4(40)に認識パケットACK14を単一伝送すると同時に、プロセッサノードPN1(10)のローカルシステムバス218を介してその書込み要求を生成するプロセッサモジュール212に該当データブロックを伝送する。また、要求したデータブロックが遠隔共有メモリに該当すると、PN1(10)は遠隔データキャッシュ224-1に該当データブロックを修正された有効な状態、例えば、「更新」状態で格納し、該当データブロックがローカル共有メモリ220に該当すると、PN1(10)はローカル共有メモリ220に該当データブロックを格納すると同時に、メモリディレクトリ226の状態を他のプロセッサノードの遠隔キャッシュが該当データブロックを共有していないことを意味する状態、例えば、「C」状態に更新する。

【0039】一方、第3番目の場合は、PN1(10)内の第1プロセッサモジュール212がデータブロックに対する書込み要求または無効化要求を生成する場合であり、図5及び図6は、この場合にもまだ有効であり、説明の便宜上、要求パケットはこの場合にもRQ12で表す。

【0040】PN1(10)が該当データブロックを遠隔キャッシュ224に有効な状態、例えば、「共有」状態で格納されている場合は、要求過程はPN1(10)が該当データブロックを遠隔キャッシュ224とローカル共有メモリ220のどこにも有効な状態で格納されていない場合と同様に行われる。

【0041】一方、図5及び図6で、プロセッサノードPN1(10)のプロセッサモジュール212及び214のいずれかからの書込み要求や無効化要求に対して、プロセッサノードPN1(10)がローカル共有メモリ220に該当データブロックを有効な状態で格納しているか、または該当ブロックを遠隔キャッシュ224に「更新-共有」状態で格納されている場合、プロセッサノードPN1(10)は第1リングバス90Aまたは第2リングバス90Bを介して他のプロセッサノードPN2~PN8(10~80)に無効化要求パケットを同報通信する。このとき、該当データブロックが奇数ブロックである場合は、図5に示すように、第1リングバス90Aを介して無効化要求が同報通信される。一方、該当データブロックが偶数ブロックである場合は、図6に示すように、第2リングバス90Bを介して無効化要求が同報通信される。要求パケットRQ12がリングバス90Aまたは90Bを巡回する間、プロセッサノードPN2~PN8(20~80)の各々は、RQ12に反応して内部の遠隔キャッシュまたはメモリディレクトリを調査して該当データブロックがどのような状態で格納されているかなどに対するスヌーピングを行うと同時に、上述したようにRQ12を隣合う次のプロセッサノードにバイパスする。前記RQ12はリングバスを巡回した後、プロセッサノードPN1(10)により除去される。一方、スヌーピングの結果、該当データブロックが他のプロセッサノード、即ちPN2~PN8(20~80)の遠隔キャッシュに更新されていない有効な状態、例えば、「共有」状態

で格納されていると表れると、PN4(40)の遠隔キャッシュの状態は、無効化された状態、例えば「無効」状態に変更する。該当データブロックが遠隔共有メモリに該当すると、PN1(10)は遠隔データキャッシュ224-1に格納された該当データブロックの状態を更新した状態、例えば、「更新」状態に変更し、該当データブロックがローカル共有メモリ220に該当すると、PN1(10)はメモリディレクトリ226の状態を他のプロセッサノードの遠隔キャッシュが該当ブロックを共有していないことを意味する状態、例えば、「C」状態に更新する。

【0042】第4番目の場合は、プロセッサノードPN1(10)の遠隔キャッシュ224におけるデータブロック置換えにより抽出されるデータブロックの状態が更新された状態、例えば、「更新」あるいは「更新-共有」状態の場合である。この場合に、PN1(10)は該当データブロックをもともと格納されるべきローカル共有メモリを備えたプロセッサノード、例えば、PN4(40)に一番近い経路を有するリングバス90Aまたは90Bを介して抽出されたブロックを含むパケットを単一通信する。そうすると、PN4(40)はRQ12に反応して内部のデータメモリとメモリディレクトリを更新し、応答パケットRSP42を第1リングバス90Aまたは第2リングバス90Bを介してPN1(10)に単一通信する。PN1(10)は認識パケットACK14をPN4(40)に単一通信する。

【0043】一方、本発明の好適な実施例によると、プロセッサノード、例えば、PN1(10)は印加されるパケットの入力順に従って1つ以上のデータブロックに対する要求を処理することができる。例えば、要求されたデータブロックに対する応答パケットが対応するプロセッサノード、即ちPN4(40)から受信される前に他のデータブロックに対する要求パケットが他のプロセッサノードPN2~PN8(20~80)中の1つ以上から伝達されると、プロセッサノードPN1(10)はまず他のプロセッサノードPN2~PN8(20~80)中の1つ以上から伝達された要求パケットに対する動作を行った数に要求したデータブロックに対する応答パケットに該当する動作を行う。

【0044】図7ないし図10を参照すると、図3に示す本発明による方向分離二重リング構造を有する多重プロセッサシステムで使用されるプロセッサノードの他の実施例がそれぞれ示されている。

【0045】図7は本発明の第2実施例によるプロセッサノード200-1の詳細図を示している。図7に示すように、プロセッサノード200-1の構成は第1リンク制御器230A及び第2リンク制御器230Bがリンクバス228無しにノード制御器222に直接接続されていること以外は、図4に示す本発明の実施例1によるプロセッサノードの構成と同様である。

【0046】図8には本発明の実施例3によるプロセッサノード200-2が詳細に示されており、プロセッサノード200-2の構成はプロセッサノード200-2がローカル共有

10

20

30

40

50

メモリ及びメモリディレクトリを含まないこと以外は、図4に示す本発明の実施例1によるプロセッサノードの構成と同様である。

【0047】図9は本発明の実施例4によるプロセッサノード200-3を詳細に示しており、プロセッサノード200-3の構成はプロセッサノード200-3が遠隔キャッシュを含まないということ以外は、図4に示す本発明の実施例1によるプロセッサノードの構成と同様である。

【0048】図10は本発明の実施例5によるプロセッサノード200-4を詳細に示しており、プロセッサノード200-4の構成は内部プロセッサモジュールがローカルシステムバスの代わりにリングまたはクロスバースイッチのようなある相互接続網240を介して互いに接続されていること以外は、図4に示す本発明の実施例1によるプロセッサノードの構成と同様である。

【0049】図7ないし図10に示す実施例の構成は、第1実施例の構成と殆ど同一であるため、その動作の説明は便宜上省略する。本発明によると、プロセッサノードPN1~PN8(10~80)の各々は図7ないし図10から選択された構成を有するプロセッサノードを介して実現することができる。

【0050】さらに、上記において、遠隔キャッシュが「更新」と「更新-共有」、「共有」、「無効」状態を有する場合について説明したが、本発明は遠隔キャッシュが変更された他の状態を有する場合を含む多様な場合にも同様に適用することができる。本発明の実施例においてローカル共有メモリのためのディレクトリが「C」、「S」、「G」の状態を維持する場合について説明したが、本発明はディレクトリが変更された多様な他の状態を有する場合にも同様に適用され得ることを理解しなければ成らない。

【0051】上述した本発明の実施例において、同報通信要求パケットに対する応答パケットと認識パケットが最短経路を有する第1リングバス90Aあるいは第2リングバス90Bを介して伝送される場合について説明したが、本発明が第1リングバス90Aを介して伝送された同報通信要求に対する応答パケットと認識パケットは第1リングバス90Aを介して伝送され、第2リングバス90Bを介して伝送された同報通信要求に対する応答パケットと認識パケットは第2リングバス90Bを介して伝送される場合と、所定の順に第1リングバス90Aあるいは第2リングバス90Bを介して応答及び認識パケットに伝送される場合にも同様に適用されることは明らかである。

【0052】

【発明の効果】従って、本発明によれば、分散共有メモリ多重プロセッサシステムはスヌーピング方式を使用することにより、プロセッサノード間のキャッシュ一貫性

を維持することができるだけでなく、リング帯域幅が2倍に拡張された単方向リングバスを有する既存システムに比べてさらに向上した性能を供給する効果を奏する。

【図面の簡単な説明】

【図1】従来技術のスヌーピング方式の単一リングを備えた分散共有メモリ多重プロセッサシステムの構成図である。

【図2】図1に示すプロセッサノードの詳細構成図である。

【図3】本発明による方向分離二重リングを備えた分散共有メモリ多重プロセッサシステムの構成図である。

【図4】本発明の実施例1による図3に示すプロセッサノードの詳細構成図である。

【図5】本発明の実施例1による分散共有多重プロセッサシステムの動作を例示的に示す図面である。

【図6】本発明の実施例1による分散共有多重プロセッサシステムの動作を例示的に示す図面である。

【図7】本発明の実施例2による図3に示すプロセッサノードの詳細構成図である。

【図8】本発明の実施例3による図3に示すプロセッサノードの詳細構成図である。

【図9】本発明の実施例4による図3に示すプロセッサノードの詳細構成図である。

【図10】本発明の実施例5による図3に示すプロセッサノードの詳細構成図である。

【符号の説明】

90A：第1リングバス

90B：第2リングバス

212：第1プロセッサモジュール

214：第2プロセッサモジュール

216：I/Oブリッジ

218：ローカルシステムバス

220：ローカル共有メモリ

222：ノード制御器

224：遠隔キャッシュ

224-1：遠隔データキャッシュ

224-2：遠隔タグキャッシュ

224-2A：第1遠隔タグキャッシュ

224-2B：第2遠隔タグキャッシュ

226：メモリディレクトリ

226A：第1メモリディレクトリ

226B：第2メモリディレクトリ

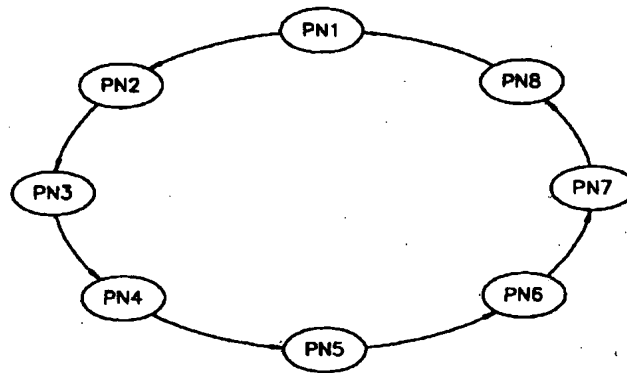
228：リンクバス

230：リングインタフェース

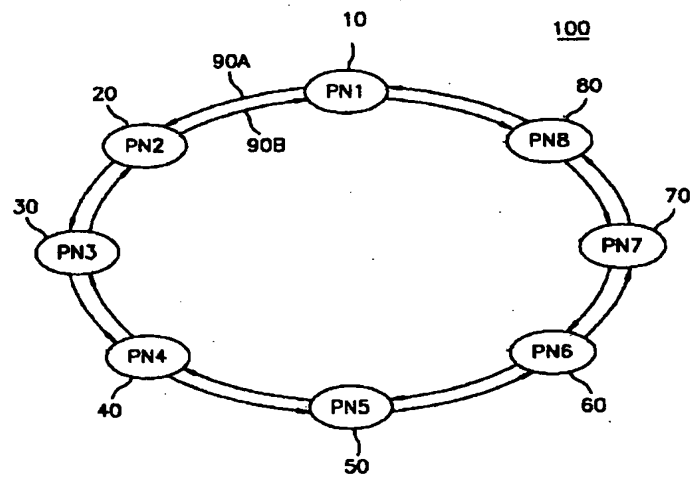
230A：第1リンク制御器

230B：第2リンク制御器

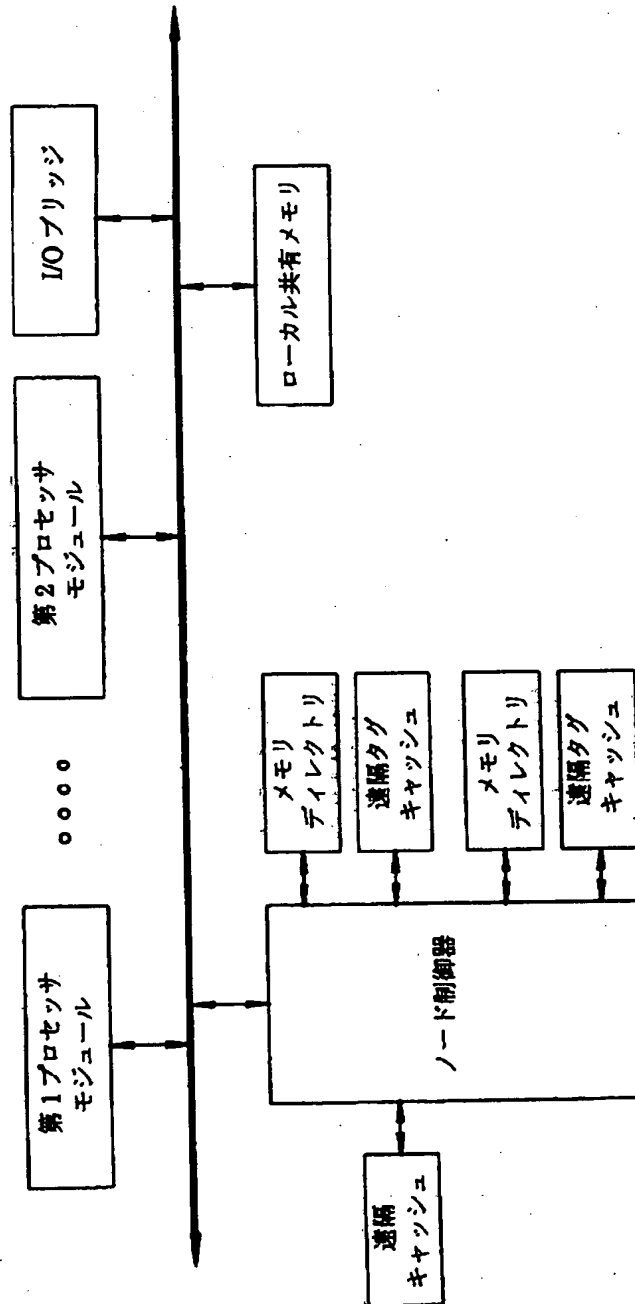
【図 1】



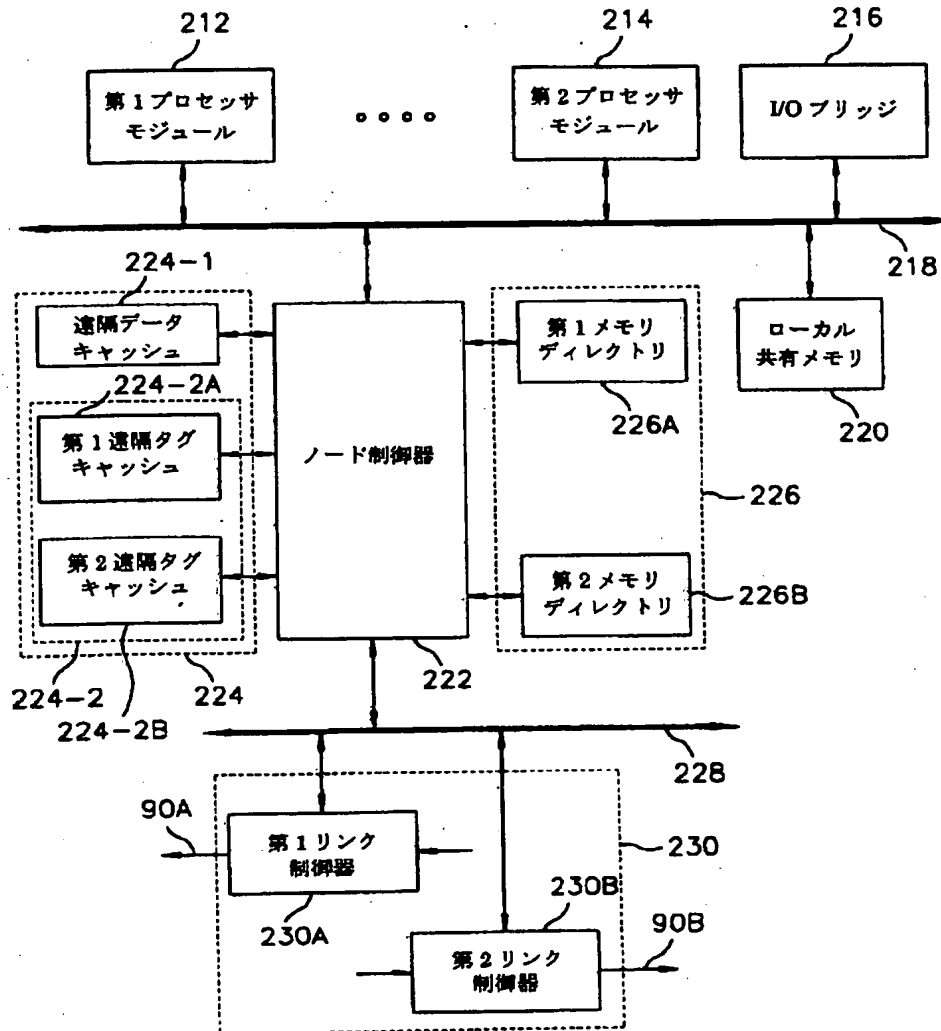
【図 3】



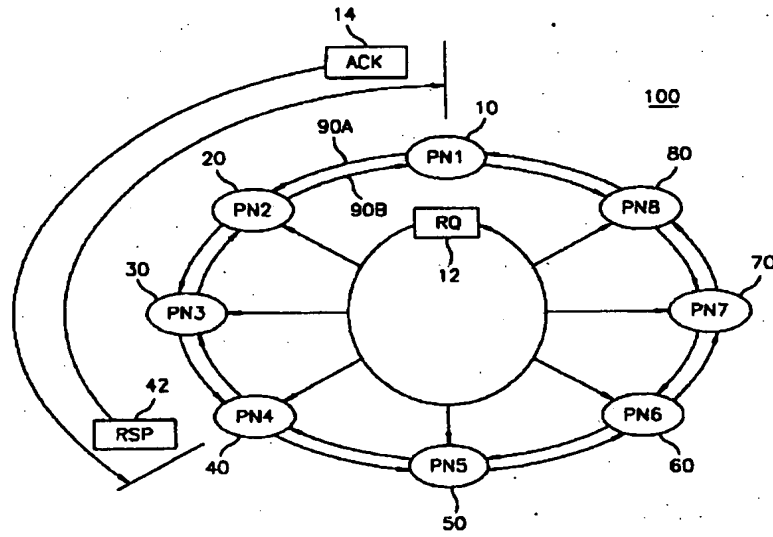
【図2】



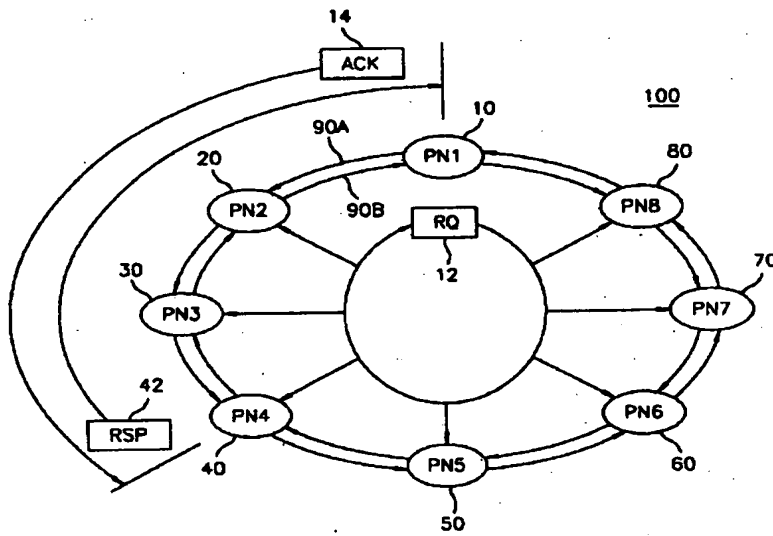
【図4】



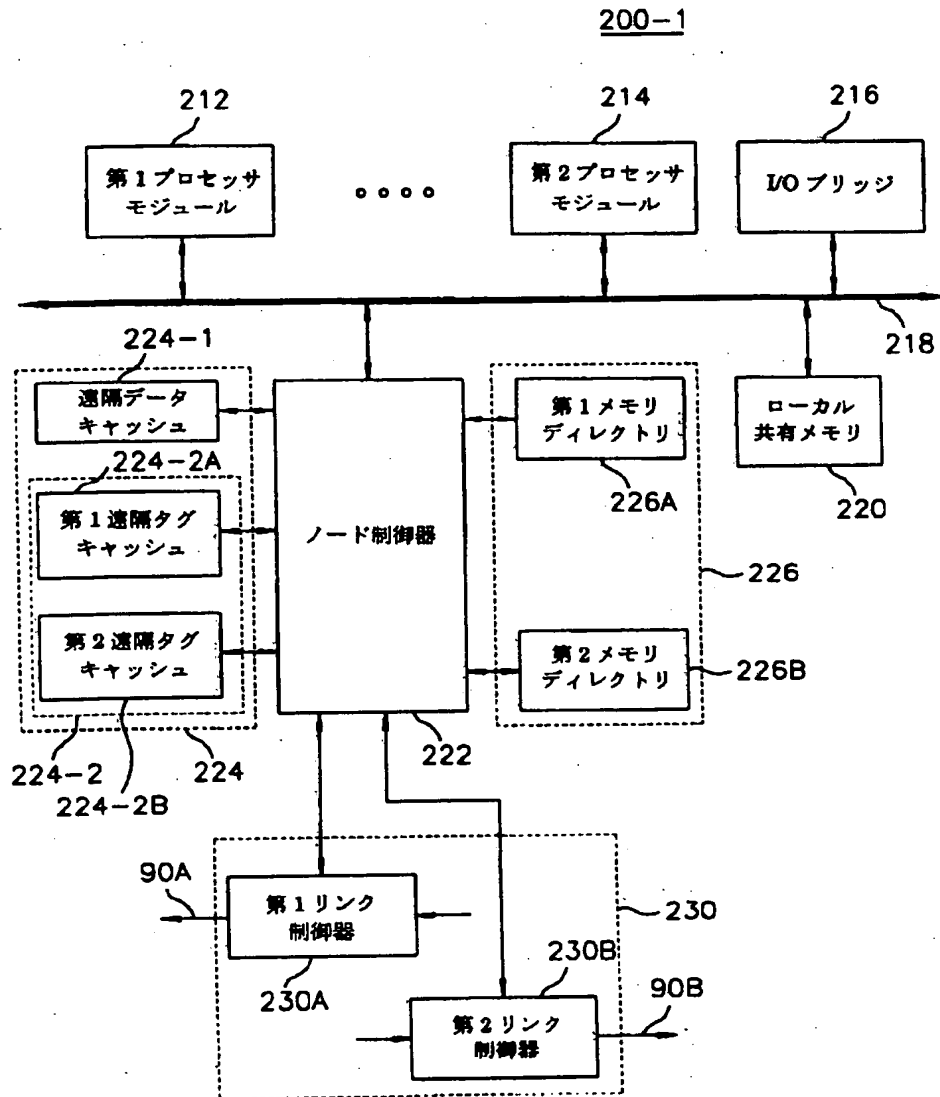
【図 5】



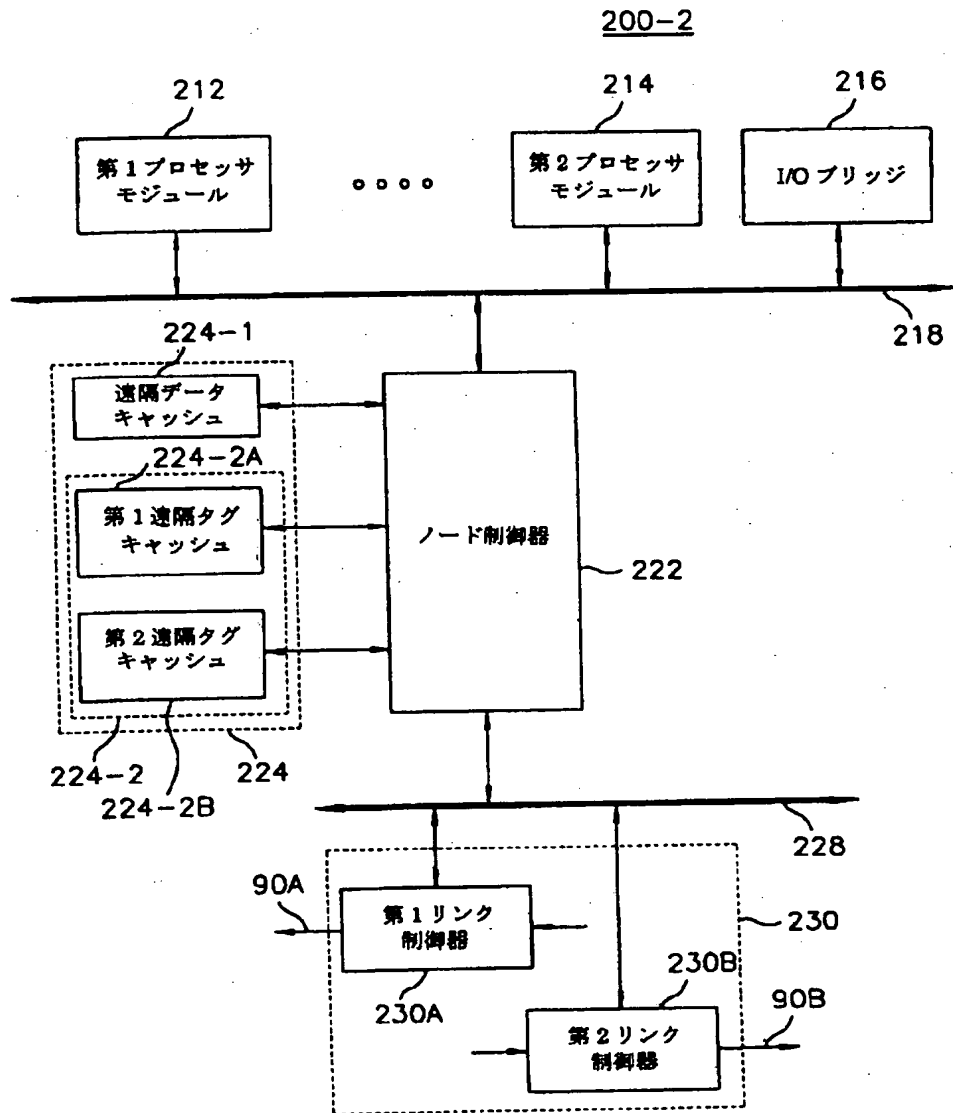
【図 6】



【図7】

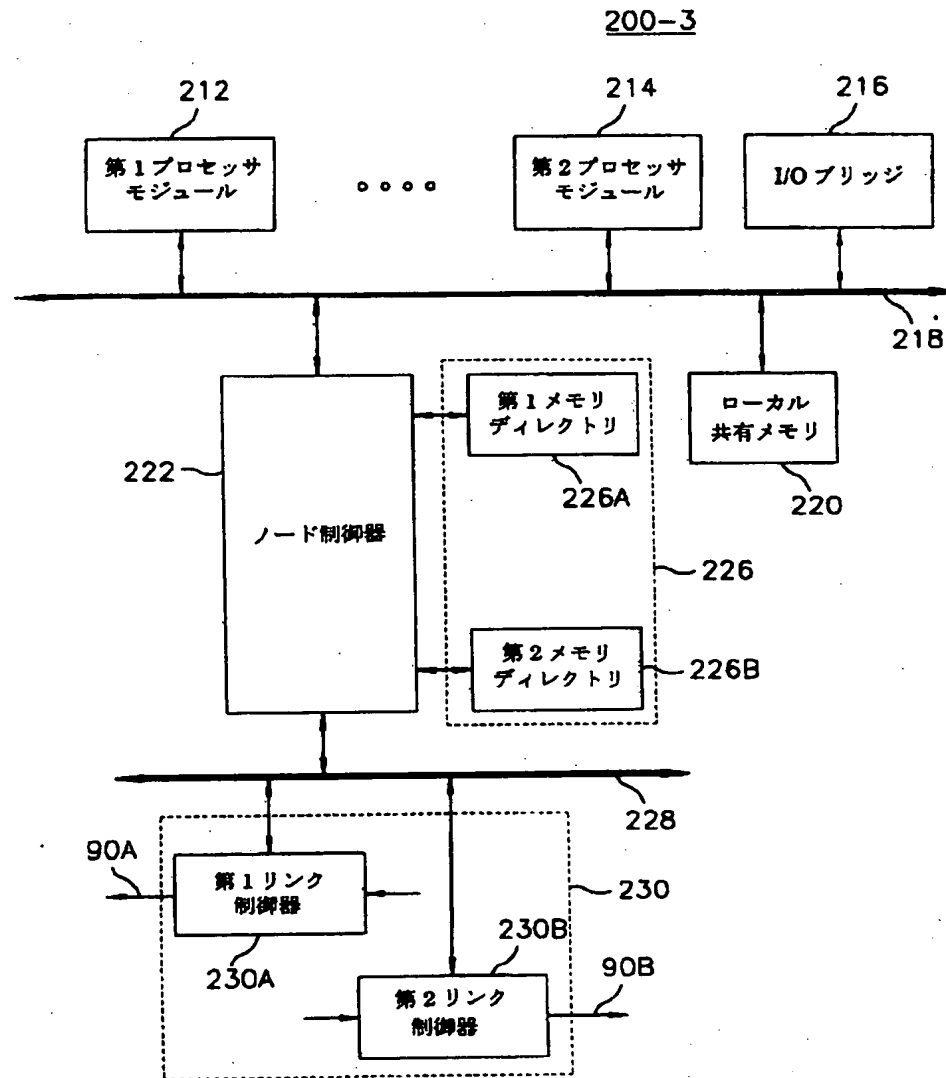


【図8】

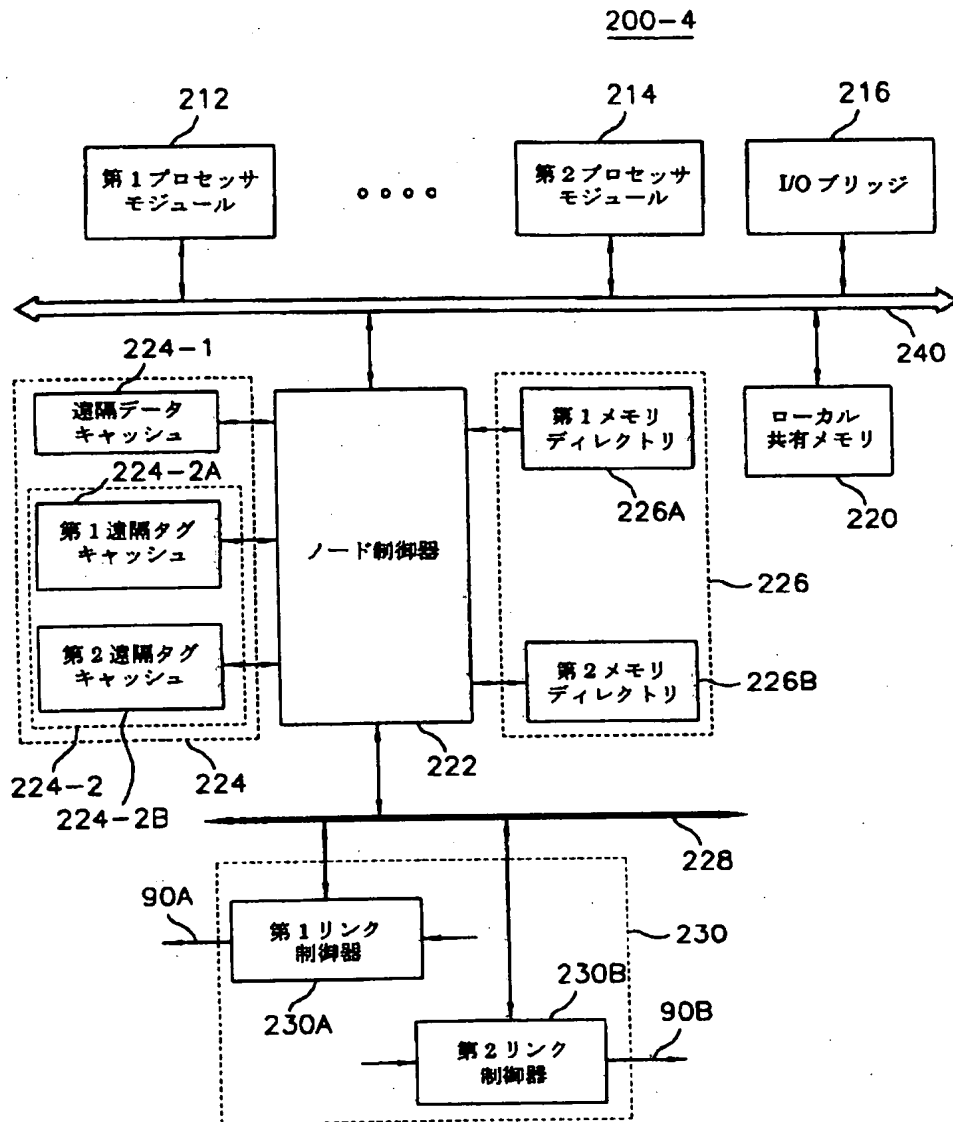




【図9】



【図10】



フロントページの続き

(71)出願人 500170526  
 金 明柱  
 大韓民国、ソウル特別市蘆原区中溪洞市営  
 アパートメント 206-304  
 (72)発明者 張 星泰  
 大韓民国、ソウル特別市城北区貞陵1洞  
 1015番地京南アパートメント103-1701

(72)発明者 全 洲植  
 大韓民国、ソウル特別市江南区道谷洞宇星  
 アパートメント1-103  
 (72)発明者 金 明柱  
 大韓民国、ソウル特別市蘆原区中溪洞市営  
 アパートメント 206-304

F ターム(参考) 5B005 JJ01 KK14 MM01 NN53 PP21  
PP26  
5B045 BB13 BB17 BB24 BB28 BB29  
BB47 DD01 DD12 GG01  
5B060 KA01 KA06